# Challenges of 3D DRAM Memories

Workshop Session on:
**3D Systems for Machine Learning and Memory architecture**

Bandwidth, Power, Temperature, Reliability

*Dr.-Ing. Christian Weis*

# Why do we care about DRAM ?
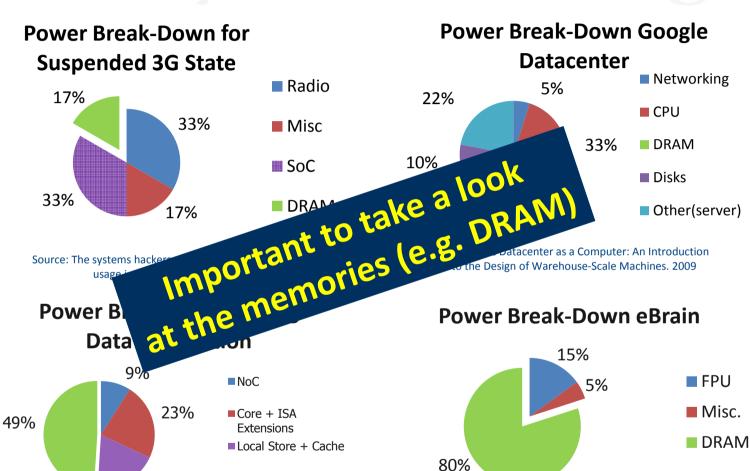
## Power Break-Down for Suspended 3G State



Legend: Radio, Misc, SoC, DRAM

17%, 33%, 33%, 17%

Source: The systems hacker... usage i...

## Power Break-Down Google Datacenter



Legend: Networking, CPU, DRAM, Disks, Other(server)

5%, 22%, 33%, 10%

...Datacenter as a Computer: An Introduction ...to the Design of Warehouse-Scale Machines. 2009

## Power B... Data... ...ion



Legend: NoC, Core + ISA Extensions, Local Store + Cache, DRAM + Controller

9%, 23%, 49%, 19%

Source: Power Consumption of Green Wave Architecture 2011
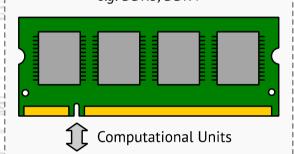
## Power Break-Down eBrain



Legend: FPU, Misc., DRAM

15%, 5%, 80%

Source: A Scalable Custom Simulation Machine for the Bayesian Confidence Propagation Neural Network model of the Brain, 2014
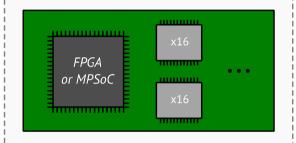
**Important to take a look at the memories (e.g. DRAM)**

MICROELECTRONIC SYSTEMS DESIGN RESEARCH GROUP

2

# Comparison of DRAM Subsystems



**DIMM Based:**

General Purpose Computers
*e.g. DDR3, DDR4*

Computational Units

**Device Based:**

Embedded / Tablets / Graphic Cards
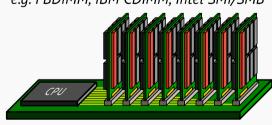*e.g. LPDDR3, GDDR5*

FPGA
or MPSoC

x16

x16

**Package on Package (PoP):**

Soldered on top of the MPSoC.
Smartphones
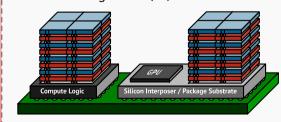*e.g. LPDDR3, LPDDR4*

DRAM

MPSoC

**Buffer on Board:**

Memory Controller on Buffer Chip,
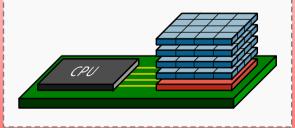Serial Connection
*e.g. FBDIMM, IBM CDIMM, Intel SMI/SMB*

CPU

**3D/2.5D-Integrated:**

Stacked on Logic or Silicon Interposer
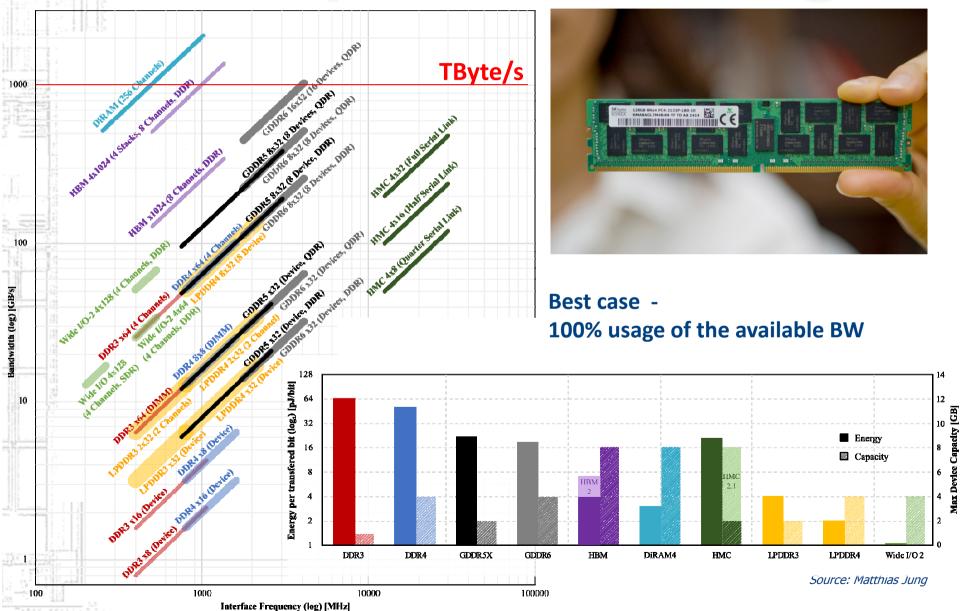by means of TSVs
*e.g. Wide I/O, HBM*

Compute Logic

GPU

Silicon Interposer / Package Substrate

**Memory Cube:**

3D-Stacked, Memory Controller on
Bottom Layer, Serial Interconnect (SerDes)
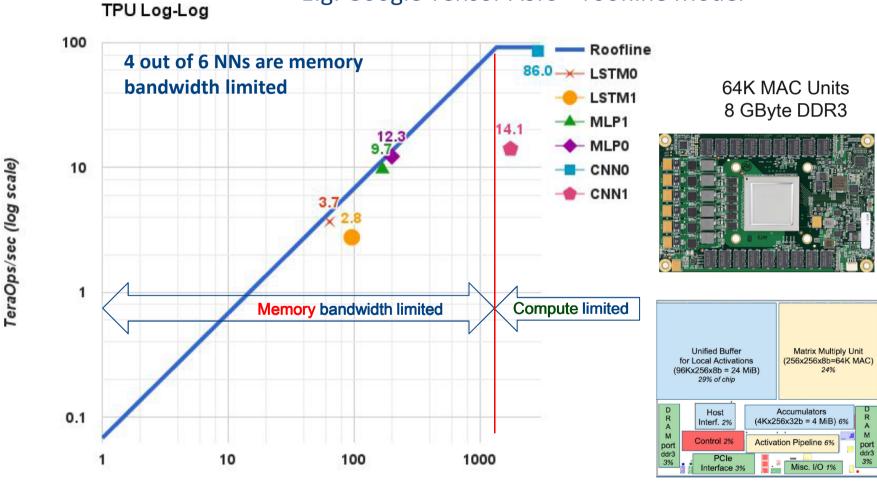*e.g. HMC, SMC*

CPU

*Source: Matthias Jung*

MICROELECTRONIC
SYSTEMS DESIGN
RESEARCH GROUP

# Comparison of DRAM Subsystems



**TByte/s**

**Best case -**
**100% usage of the available BW**

*Source: Matthias Jung*

# Bandwidth Challenge

### E.g. Google Tensor ASIC – roofline Model

**TPU Log-Log**

**4 out of 6 NNs are memory bandwidth limited**

Roofline
LSTM0 — 86.0
LSTM1
MLP1 — 9.7
MLP0 — 12.3
CNN0
CNN1 — 14.1

3.7 — 2.8

TeraOps/sec (log scale)

Memory bandwidth limited — Compute limited

Operational Intensity: Ops/weight byte (log scale)

64K MAC Units
8 GByte DDR3

Unified Buffer for Local Activations (96Kx256x8b = 24 MiB) 29% of chip

Matrix Multiply Unit (256x256x8b=64K MAC) 24%

DRAM port ddr3 3%

Host Interf. 2%

Accumulators (4Kx256x32b = 4 MiB) 6%

DRAM port ddr3 3%

Control 2%

Activation Pipeline 6%

PCIe Interface 3%

Misc. I/O 1%

# Bandwidth Challenge in 3D-DRAMs

HBM2 and HMC can provide huge bandwidths.

BUT, there is a price to pay …



**>256 GB/s**
**Bandwidth / Cube**

**High Bandwidth Memory (HBM2)**

**Hybrid Memory Cube (HMC)**

**>240 GB/s**
**Bandwidth / Cube**

**Power: 40-50 GB/s/W / Cube**

**Power: ≈20 W / Cube**

# HMC Power >>11W

**DRAM part only power:**

**for different page sizes and technologies**

**Link-power only is about 10-11W!**

# Detailed DRAM Energy Distribution

- DRAM Power Breakdown for Twitter Memcached Application*
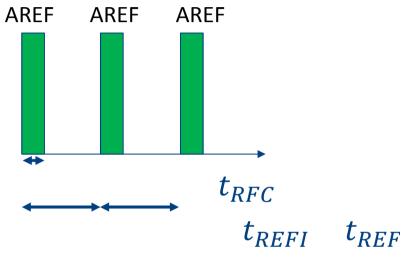- 2GB LPDDR3 (Low-Power DDR3 DRAM)

- **Refresh optimizations**
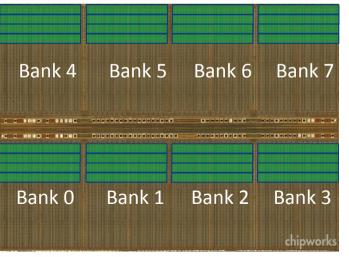- **Minimizing Row misses**



Pie chart segments: Refresh, Activate, Precharge, Read, Write

- **Using 3D Interconnect**
- **Maximize useful data**

B0  B2  B4  B6

PSAs

SSAs (e.g. 8B)

e.g. 1kB

x8

B1  B3  B5  B7

## Important DRAM Commands

| | | |
|---|---|---|
| ACT: Activates a specific row in a specific bank (sensing into PSA) | tRCD |
| RD: Read from activated row (prefetch from PSA to SSA and burst out) | tCL + tBURST |
| PRE: Precharges set LWL=0 set LBL=VDD/2 | tRP |
| REFA: DRAM cells are leaky and have to be refreshed | tREFI & tREF |

# How Refresh is Performed?

- DRAM controller sends AREF commands every $t_{REFI}$ (eg. 7.8 µs for Temp. < 85°C)

- Single AREF command refreshes multiple rows in all banks ( eg. '2' rows in all 8 banks for 2Gb DDR3 DRAM)
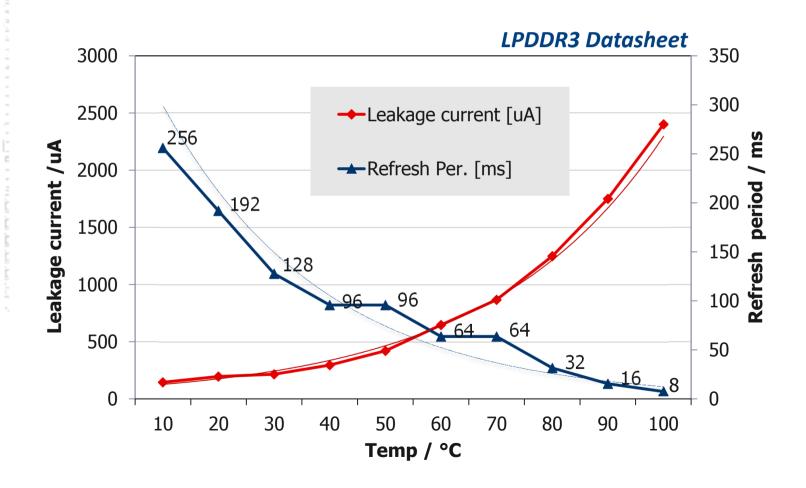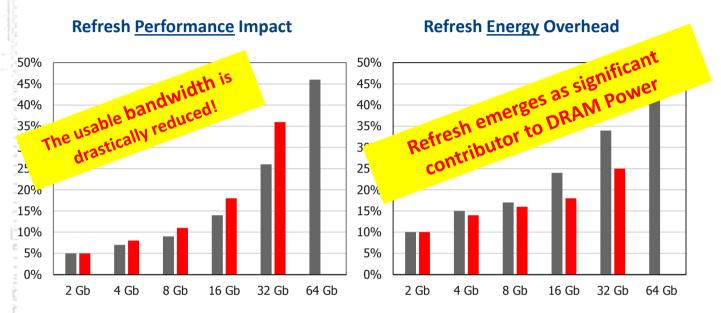
AREF    AREF    AREF

$t_{RFC}$

$t_{REFI}$    $t_{REF}$

| Bank 4 | Bank 5 | Bank 6 | Bank 7 |

| Bank 0 | Bank 1 | Bank 2 | Bank 3 |

chipworks

# Refresh/Temperature Challenge

**Exponential temperature/leakage current behavior**
**→ shorter refresh periods**



*LPDDR3 Datasheet*

# Refresh/Temperature Challenge

**Refresh Performance Impact**

The usable bandwidth is drastically reduced!

**Refresh Energy Overhead**

Refresh emerges as significant contributor to DRAM Power

■ *J. Liu, et al. RAIDR: Retention-Aware Intelligent DRAM Refresh, ISCA 2012*
■ *I. Bhati, et al. DRAM Refresh Mechanisms, Trade-offs and Penalties, IEEE Trans. 2015*

4 TB DDR3 DRAM
Stand-by 300W (only Refresh)

*Paul Rosenfeld (IBM Server on display at Supercomputing)*

# Refresh Optimization Techniques

To counterbalance this trend for future devices and the higher temperatures in 3D-DRAMs …

- **Temperature-aware Bank-wise** Refresh (detailed control)
- **Approximate** DRAM
- **ORGR** – Optimized Row Granular Refresh  (only refresh data that is stored in an optimized way)
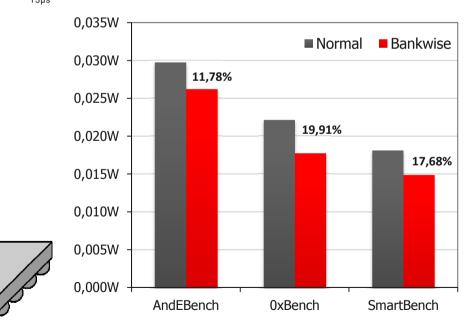
# Temperature-aware Bank-Wise Refresh



- Different refresh rates on different dies (bank groups), according to the temperature of the die/bank
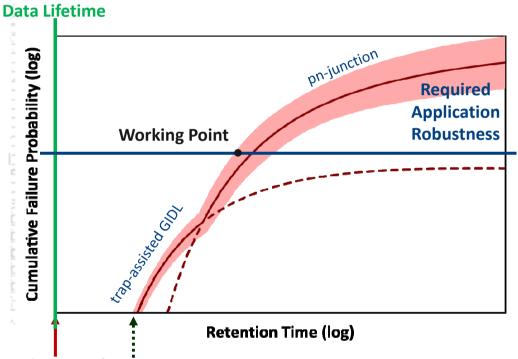
- Each bank was equipped with a TS (Temp Sensor)

*M. Sadri, et al. Energy Optimization in 3D MPSoCs with Wide-I/O DRAM Using Temperature Variation Aware Bank-wise Refresh, DATE 2014*

13

# Approximate DRAM

Lowering or completely switching off refresh, accepting risk of data errors

- Consider DRAM device as a stochastic model that includes process variations

**Data Lifetime**

*Y-axis:* Cumulative Failure Probability (log)

*X-axis:* Retention Time (log)

pn-junction

**Working Point**

trap-assisted GIDL

**Required Application Robustness**

**Conservative Datasheet Refresh Period Guardband (i.e. 64ms)**

**Required Refresh Period based on measurements (First errors happen after e.g. 1s)**

## Switch Off Refresh

- If data lifetime is smaller than required refresh period
- If data lifetime is larger than required refresh period AND application has some robustness w.r.t. errors

*Statistical retention error model and measurements mandatory*

MICROELECTRONIC SYSTEMS DESIGN RESEARCH GROUP

# DRAMMeasure DDR4
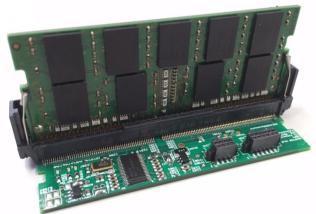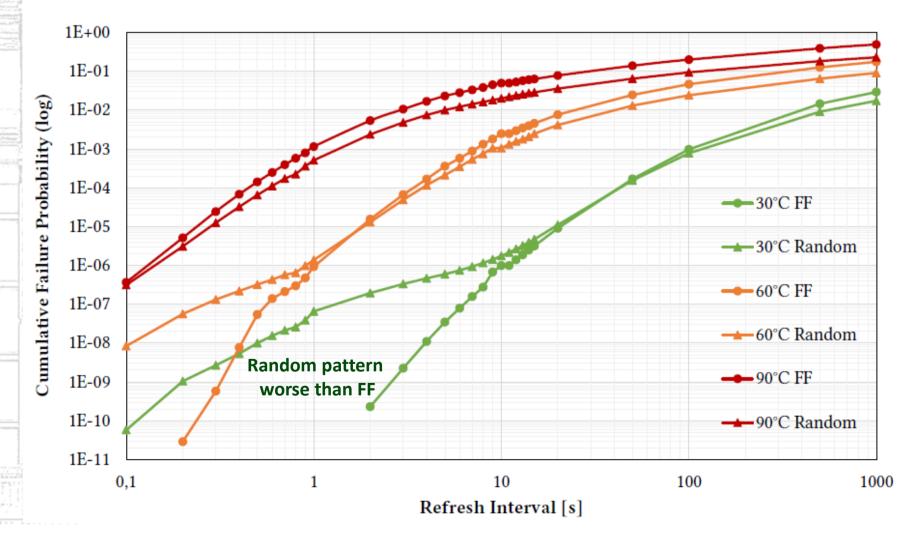


- High freq. 1.2GHz and higher: DDR4-2400
- Precise temperatures for heating up to 95°C
- Exact current measurements incl. VPP

- Retention behavior depends on cell leakage (drain, sub-threshold, cell capacitor), cross talk, process variations, temperature, cell type
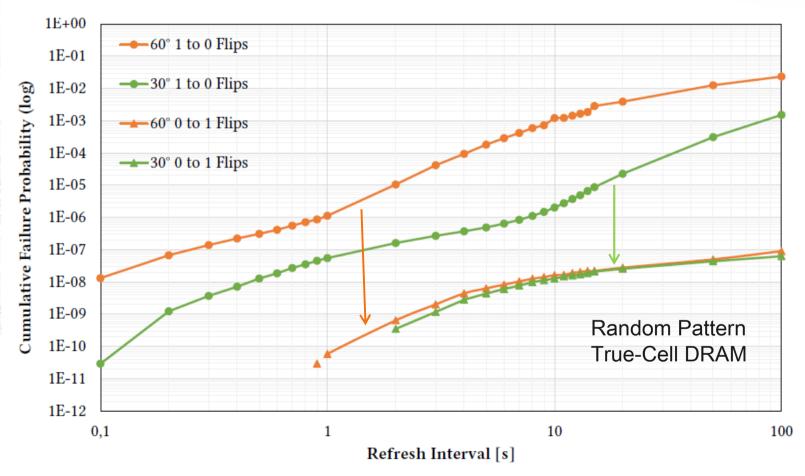
# DDR4 Retention Time Measurements II

- Unsymmetrical error behavior dependent on cell type  (true-cell, anti-cell)



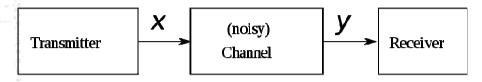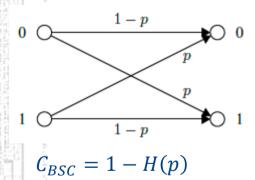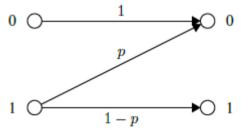**1→0 flip much more likely than 0→1 flip**

Consider memory as noisy communication channel



Symmetric retention behavior:
*Binary Symmetric Channel (BSC)*



$$C_{BSC} = 1 - H(p)$$

Asymmetric retention behavior:
*Z-Channel*



$$C_Z = \log_2\left(1 + (1-p) \cdot p^{\frac{p}{1-p}}\right) \approx 1 - \frac{1}{2}H(p)$$

- Larger reliability if internal cell structure is known
- More efficient **ECC** techniques possible
- Appropriate data representation: e.g. small dynamic range C2 versus sign/magnitude
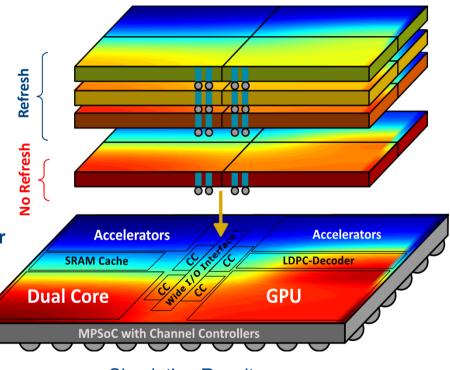
# A Per Layer Refresh Policy for 3D DRAMs

Separation of a 3D DRAM Stack into **unreliable** and **reliable** regions

- Reliable regions: higher DRAM layers with temperature aware refresh
- Unreliable region: bottom DRAM layer with disabled refresh → **Omit Refresh (OR)**
- Access unreliable region while reliable region is refreshed

**Typical example applications**

- Graph processing
- Image processing
- Baseband processing

→ **Saves 100% refresh power in the unr.-layer**
→ **Increases bandwidth**



Simulation Results

*Matthias Jung, et al. 2015. Omitting Refresh: A Case Study for Commodity and Wide I/O DRAMs.*

# Reducing Refresh Overhead

- **Selective Refresh**
- **Retention Aware Refresh**
- **Approximate DRAM**

> **Different rows need to be refreshed at different rates**

Drawbacks of normal Auto-Refresh:

- AREF lacks flexibility
- No access to internal refresh row counter
- No rows can be skipped
- The complete DRAM has to be refreshed in the same rate
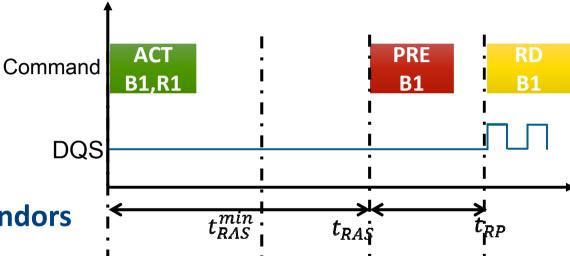
Use

**O**ptimized

**R**ow

**G**ranular

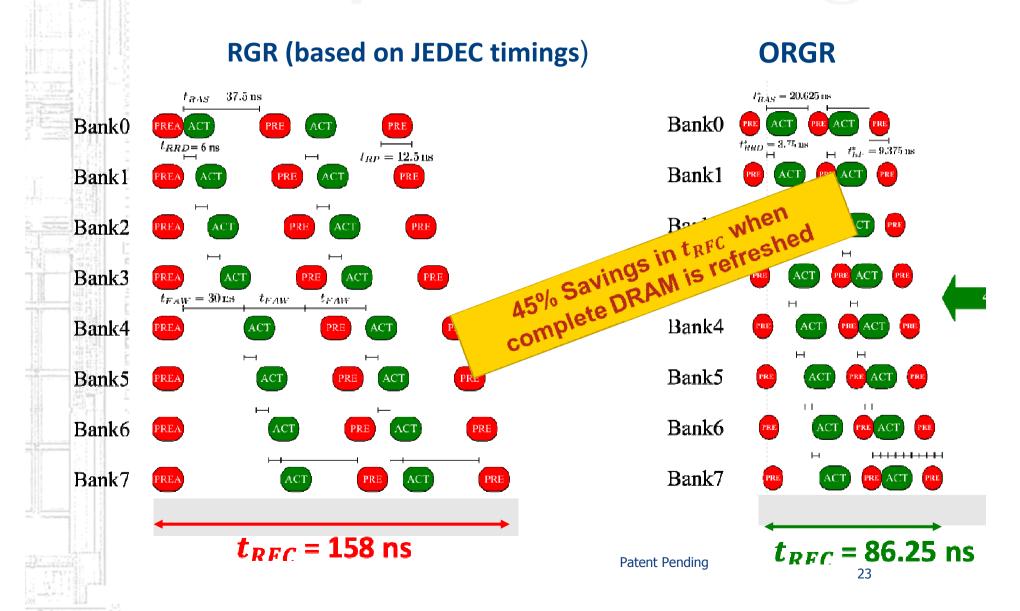**R**efresh

# ORGR – Idea / Vendor Specific



- $t_{RASmin}$ **timer at all vendors present (safety)!**

- **But, vendor specific implementation**

- **Reverse Engineering technique performed during init, boot or run-time**

$$t_{RAS} = 37.5 \text{ ns}$$
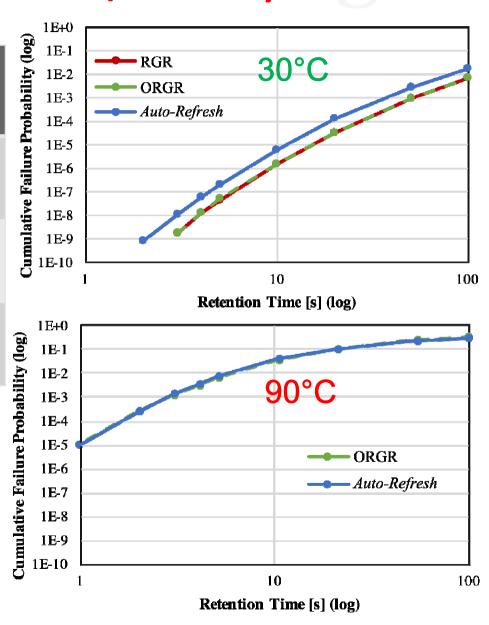$$t_{RAS}^{min} = 20.7 \text{ ns}$$

Patent Pending

22

# ORGR - Benefits



**RGR (based on JEDEC timings)**

**ORGR**

45% Savings in $t_{RFC}$ when complete DRAM is refreshed

$t_{RFC}$ = 158 ns

$t_{RFC}$ = 86.25 ns

Patent Pending

23

# ORGR – Validation/Reliability

| Refresh Technique | $t_{RFC}$ /ns | Refresh Energy /mJ |
|---|---|---|
| **Auto Refresh** | **262.5** | **186.24** |
| **RGR** | **292.5** | **230.48** |
| **ORGR** | **146.25** | **209.72** |

- Measured for 4Gb x16 DDR3 DRAM
- Refreshing the complete DRAM

# "Non-deterministic" DRAM Timing Behavior

**Chstone ADPCM** / DDR3 / **BRC** / FCFS

**Chstone ADPCM** / DDR3 / **RBC** /



*Up to 10x variance in access times*

*Row misses in the same bank*

- DRAM latency varies largely
- Depending on
  - **Application**
  - **Address Mapping**
  - **DRAM**
  - **Memory Controller**
  - ...

**Chstone ADPCM** / **WIDE IO** / BRC /

**Chstone GSM** / DDR3 / RBC / FCFS

- Similar variation in energy/DRAM access

# Application Aware Address Mapping

## Standard Mapping (BRC)



## Bank Interleaving (RBC)



## Bit Reversal Address Mapping



## Permutation-Based Page Interleaving



## Toggling Rate Analysis



Minimize row misses in the same bank



Input Address → Any Bijective Mapping ← Application Knowledge → Output Address
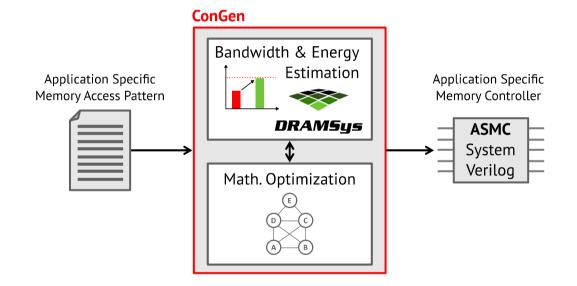
# ConGen Methodology

Exploit full application knowledge i.e. determinism of access pattern

- Minimize #row misses in same bank

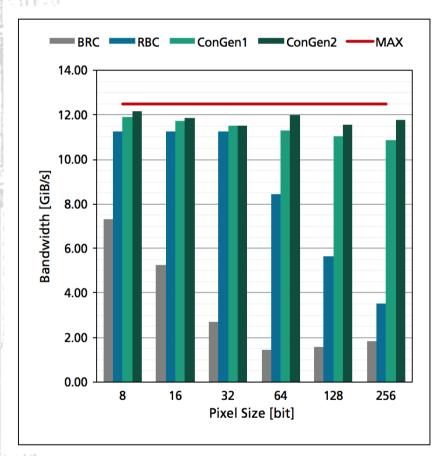- Decrease energy and latency, and increases bandwidth
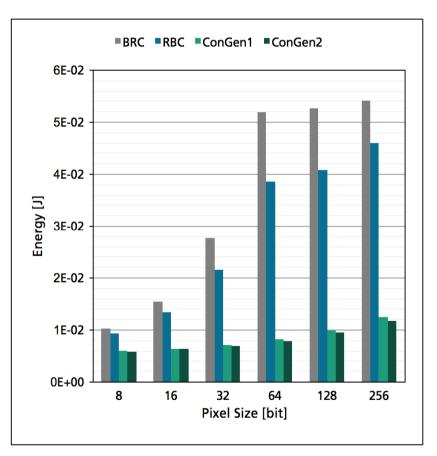


Optimization Problem

- Minimize number of row misses for all DRAM banks over an given logical memory access trace

- NP-hard problem

# ConGen Methodology - Results

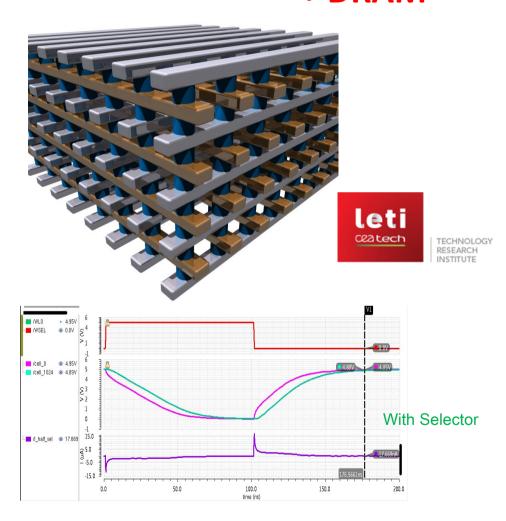## Industrial image processing task (Image Rotation 1024x 576 Pixel)



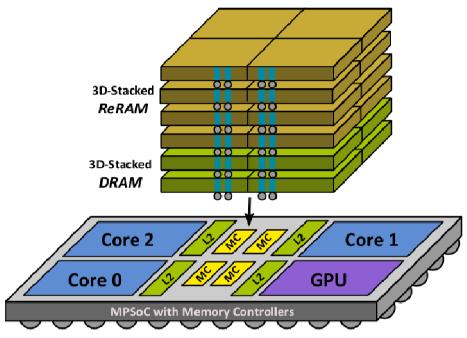Legend: BRC = Bank-Row-Column, RBC = Row-Bank-Column address mapping

# Heterogeneous Memory using ReRAM + DRAM

- Architecture level design space exploration tool to estimate the Timings, Area, and Power for high density ReRAM crossbar devices (**ReRAMSpec**)

- System Level (**SystemC**) and behavioral (**SystemVerilog**) modelling of ReRAM devices

- Circuit level modelling with **SPICE** of the ReRAM Array cross-section and periphery (Drivers, Sense-Amps etc.)
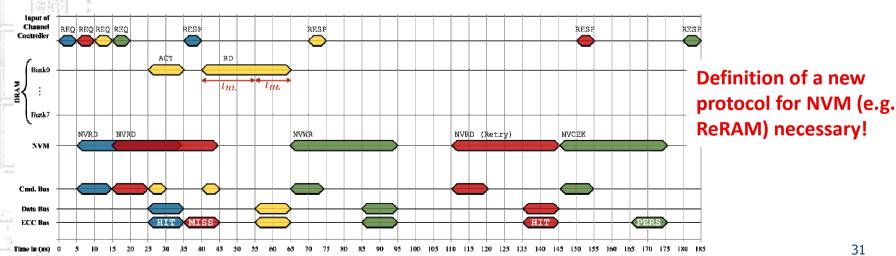


With Selector

# Heterogeneous 3D Memory System



- **Each channel consists of DRAM (smaller capacity) and ReRAM (larger capacity)**
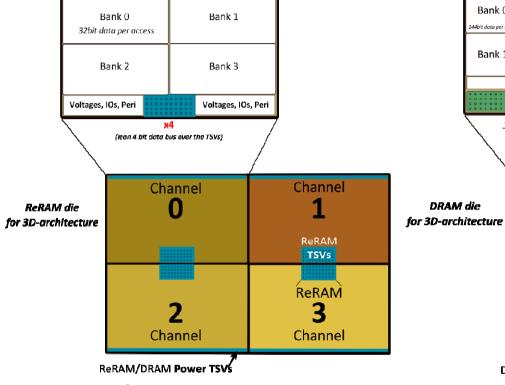- **Special Memory Controllers (MCs) needed / Hybrid**

**Two Options:**

- **DRAM as a Cache for ReRAM**
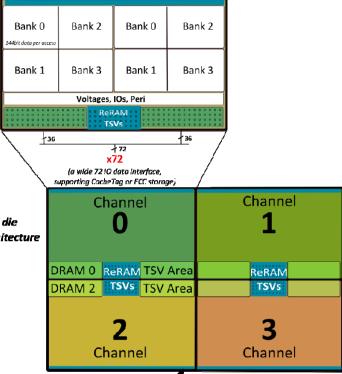- **DRAM and ReRAM individually addressable**

**Definition of a new protocol for NVM (e.g. ReRAM) necessary!**

# Heterogeneous 3D Memory System



ReRAM die for 3D-architecture

DRAM die for 3D-architecture

### ReRAM Parameters for a Heterogeneous 3D Memory System

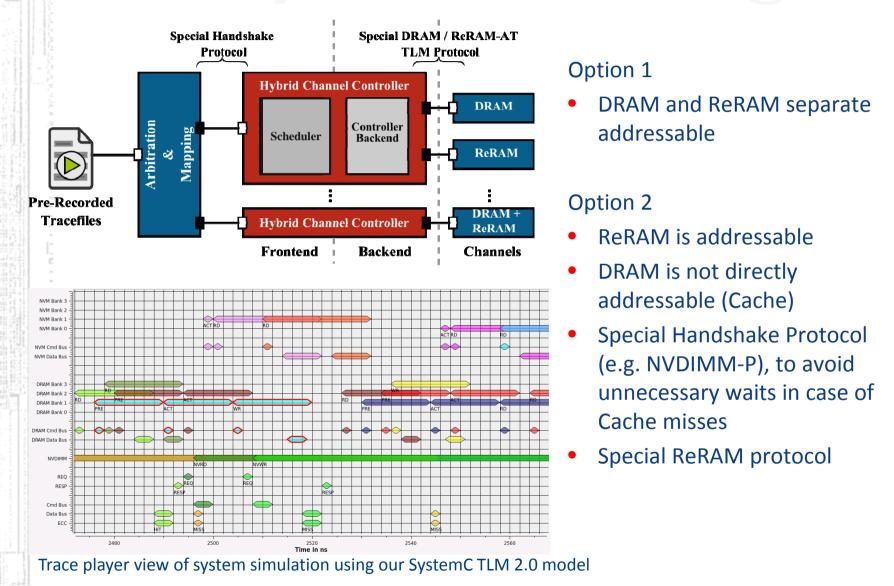| | |
|---|---|
| Capacity of the Die | 16 Gbit |
| Die Size | 11.7 x 12.6 mm |
| Number of Channels | 4 |
| Capacity/Channel and Layer (tier) size | 4 Gbit |
| ReRAM/Logic Technology | 28 nm |
| Number of Banks per Channel | 4 |
| IO width(s) | 4 data lines |
| Interface and Frequency | DDR, Prefetch 8, 500 MHz |
| Maximum Bandwidth | 1Gb/s per Pin |

### DRAM Parameters for a Heterogeneous 3D Memory System

| | |
|---|---|
| Capacity of the Die | 8 Gbit + 1 Gbit (for Tag or ECC) |
| Die Size | 9 x 11.5 mm |
| Number of Channels | 4 |
| Capacity/Channel + Layer size | 2 Gbit + 256 Mbit (for Tag or ECC) |
| DRAM Technology | 22 nm |
| Number of Banks per Channel | 4 |
| Page size(s) | 2K Bytes |
| IO width(s) | 72 data lines |
| Interface and Frequency | DDR, Prefetch 4, 500 MHz |
| Maximum Bandwidth | 1 Gb/s per Pin |

# Heterogeneous 3D Memory System



Trace player view of system simulation using our SystemC TLM 2.0 model

**Option 1**

- DRAM and ReRAM separate addressable

**Option 2**

- ReRAM is addressable
- DRAM is not directly addressable (Cache)
- Special Handshake Protocol (e.g. NVDIMM-P), to avoid unnecessary waits in case of Cache misses
- Special ReRAM protocol
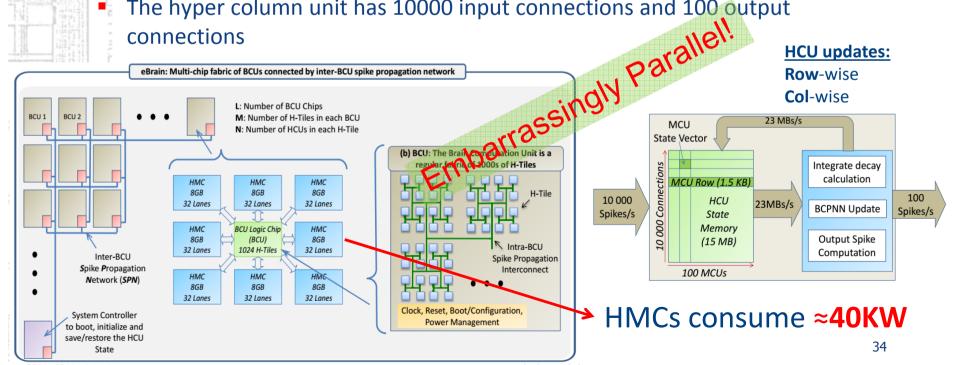
# Custom 3D-DRAM for eBRAIN II

- A custom multi-chip design to simulate the human brain in real time using the spiking BCPNN (Bayesian Confidence Neural Network )

- The architecture for this algorithm is based on Hyper Columns Units (HCU) and Mini Columns units (MCU)

- The parallel computability of HCUs and MCUs makes this architecture hardware friendly

- Each HCU is an aggregation of 100 MCUs

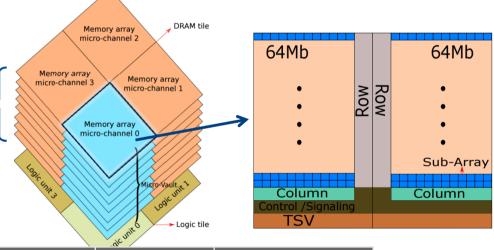- The hyper column unit has 10000 input connections and 100 output connections



eBrain: Multi-chip fabric of BCUs connected by inter-BCU spike propagation network

L: Number of BCU Chips
M: Number of H-Tiles in each BCU
N: Number of HCUs in each H-Tile

Embarrassingly Parallel!

**HCU updates:**
**Row**-wise
**Col**-wise

HMCs consume ≈**40KW**

# Custom 3D-DRAM for eBRAIN II

- Custom-optimized **3D-DRAM architecture** => 48 I/O DDR microChannel per HCU (1 – 2 mm$^2$ depending on the DRAM tech.) with 500MHz freq.

- Tailored **access** → using a technique called "**Row merge**", where we balanced the BW between Row-updates and Col-updates.
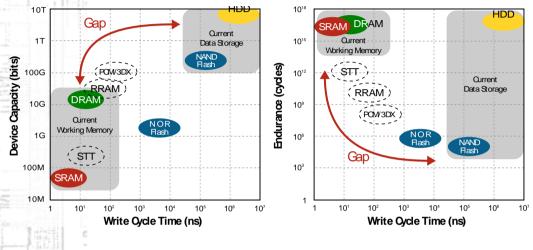


| Species | # of HCUs | Average Power |
|---------|-----------|---------------|
| Mouse | $1.6 \times 10^3$ | 13 W |
| Rat | $5.0 \times 10^3$ | 44 W |
| Cat | $6.0 \times 10^4$ | 522 W |
| Macaque | $2.0 \times 10^5$ | 1700 W |
| Human | $2.0 \times 10^6$ | **17 KW** |

Matrix – Bank mapping of 4 HCUs:
→ *optimized data layout*

# The Future is Hybrid/Heterogeneous
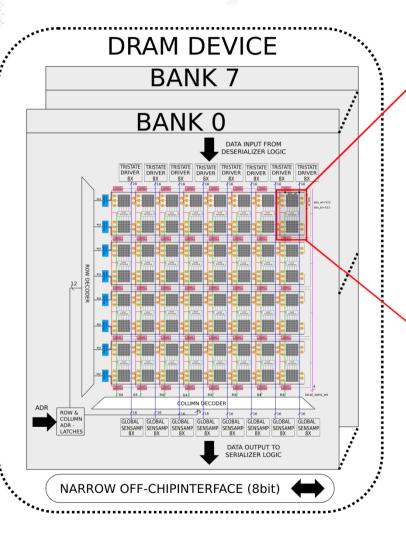


- **New memory technologies:**
  - PCM
  - 3DXPoint
  - STT-MRAM
  - **ReRAM**
- DRAM **won't** be dead, but will change its role → maybe used as **Cache ...**
- New memory **ECC** techniques
- Heterogeneous main memory systems:
  - **NVDIMM-P**
  - 3D MPSoCs / **3D Memory Stacks**
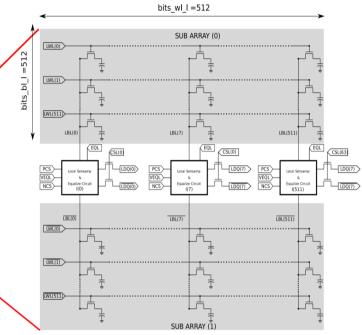- New requirements on:
  - Compiler
  - OS
- **Processing in memory (PIM)**
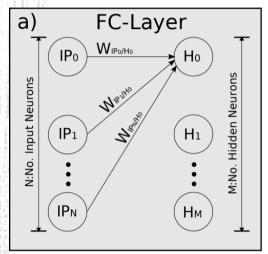
# In/Near–Memory Processing



- For NN processing (e.g. Mult & Add)
- Place special Logic between the sub-arrays
- Maximize degree of parallel data processing e.g. 512/1024 bit in DDR3/4 devices
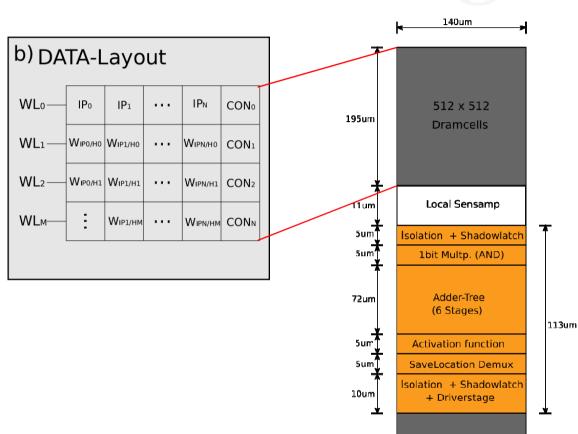
# In/Near–Memory Processing



a) **FC-Layer**

IP$_x$: Input Layer
W$_{x/y}$: Weight
H$_y$: Hidden/Output Layer
CON$_x$: Config Flag Layer Type
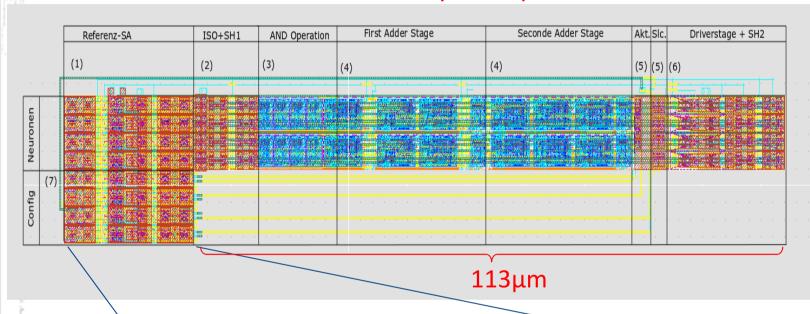WL$_x$: Wordline (Mem. Row)

b) **DATA-Layout**

**The data layout:**

- Weights correspondent to separate neurons (H0, H1, ..) are stored row-wise.
- Process the inputs of a single neuron in parallel by reading the weights in parallel from the corresponding DRAM row
- Currently 100 values (100 bits because of 1-bit weights) are accessed per row per clock cycle
- Maximum bandwidth per clock cycle can be 512/1024 bits for DDR3/DDR4 minus the control bits!.

# Layout Study – Feasibility
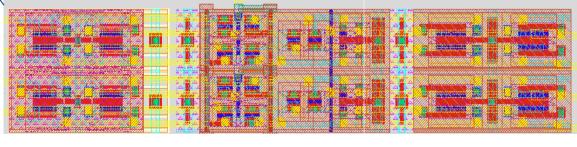
## Neural Network Layout Implementation



| Referenz-SA | ISO+SH1 | AND Operation | First Adder Stage | Seconde Adder Stage | Akt. | Slc. | Driverstage + SH2 |
|---|---|---|---|---|---|---|---|

113µm

Reference Sense-Amplifier

11µm

# In/Near–Memory Processing

Current implementation:

| | |
|---|---|
| Input data precision, QI [bits] | 2 |
| Weights precision, QW [bits] | 1 |
| Number of neurons per fully-connected layer, N | 100 |
| Number of layers, L | 10 |
| Access time per DRAM row, T [ns] | **20** |
| Number of sub-arrays per bank, S | 16 |
| Number of banks in DRAM device, B | 8 |

Taking into account that the current implementation allows to access N weights per sub block in S sub-arrays in B banks of a single device in parallel, the computed throughput is:

$$N * 2^1 * S * B / T = \textbf{1.28 TOP/s}$$

[1] 2 comes from addition and multiplication considered as separate operations

# Summary – Take-away messages

- Approximate DRAM and optimized Refresh control can be used to trade-off BW vs. reliability

- ConGen methodology to improve BW and energy

- Custom 3D-DRAMs have a large potential

- Hybrid/Heterogeneous architectures and Near/In-memory processing will be key

*Thank you for Listening*

*For more information //ems.eit.uni-kl.de*

*For the tools*

*https://www.uni-kl.de/3d-dram/tools/*

41